

# Information Extraction

–

## Architecture and Task Definition

PD Dr. Günter Neumann

DFKI and Saarland University



# Next Lecture at 10th Sept

- Geb. C7.2 Seminarraum



# Information Extraction (IE)

The goal of IE research is to build systems that find and link *relevant* information from NL text ignoring irrelevant information.

## Core Functionality

### Input

Templates coding relevant information, e.g. company, product, medical information  
set of real world texts

### Output

set of instantiated templates filled with relevant text fragments  
(normalized to a canonical form)



# Example: Job Advertisement

## Input

### Posting from Newsgroup

Telecommunications. SOLARIS Systems Administrator. 38-44K. Immediate need

Leading telecommunications firm in need of an energetic individual to fill the following position in the Atlanta office:

SOLARIS SYSTEMS ADMINISTRATOR  
Salary: 38-44K with full benefits  
Location: Atlanta Georgia, no  
relocation assistance provided

## Output

### Filled Template

computer\_science\_job  
title: SOLARIS Systems Administrator  
salary: 38-44K  
state: Georgia  
city: Atlanta  
platform: SOLARIS  
area: telecommunications



# Example: Terrorists actions

„Salvadoran President-elect Afredo Cristiani condemned the terrorist killing of Attorney General Roberto Garcia Alvarado and accused the Farabundo Marti National Liberation Front (FMLN) of crime.“

Incident: KILLING  
Perpetrator: „terrorist“  
Confidence: —  
Human Target: „Roberto Garcia Alvarado“

Incident: KILLING  
Perpetrator: FMLN  
Confidence: Accused by Authorities  
Human Target: —

Incident: KILLING  
Perpetrator: FMLN  
Confidence: Accused by Authorities  
Human Target: „Roberto Garcia Alvarado“



# Example: Company's turnover

Lübeck (dpa) – Die **Lübecker Possehl-Gruppe**, ein im Produktions-, Handel- und Dienstleistungsbereich tätiger Mischkonzern, **hat 1994** den **Umsatz** kräftig um **17 Prozent** auf rund **2,8 Milliarden DM gesteigert**. In das neue Geschäftsjahr sei man ebenfalls „mit Schwung“ gestartet. Im **1. Halbjahr 1995** hätten sich die **Umsätze** des Konzerns im Vergleich zur Vorjahresperiode um **fast 23 Prozent** auf rund **1,3 Milliarden erhöht**.

Type:	turnover
C-name:	Possehl1
Year:	1994
Amount:	2.8e+9DM
Tendency:	+
Diff:	+17%

Type:	turnover
C-name:	Possehl1
Year:	1995/1
Amount:	1.3e+9DM
Tendency:	+
Diff:	+23%



# Text-based Information Extraction (IE)

Template:

ManagementSuccession

PersonIn: \_\_\_\_\_

PersonOut: \_\_\_\_\_

Position: \_\_\_\_\_

Organisation: \_\_\_\_\_

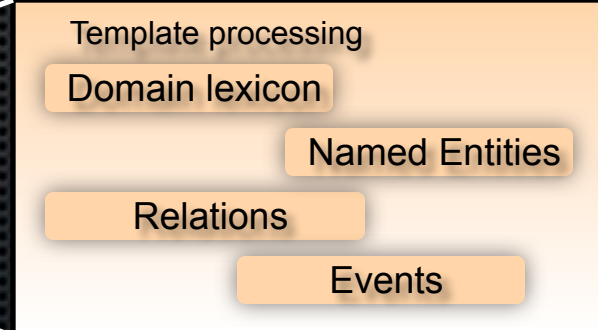
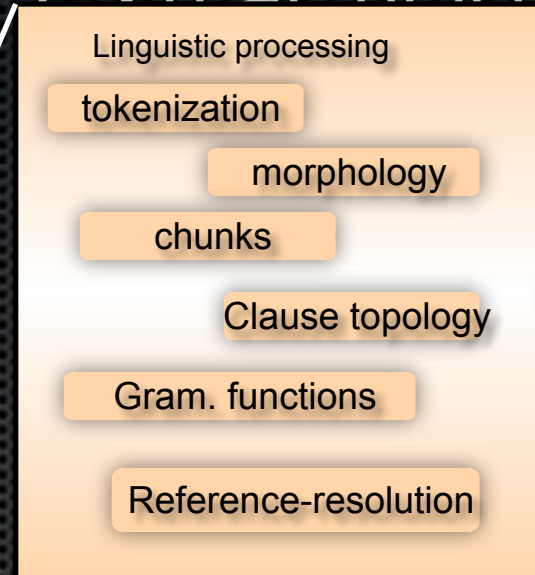
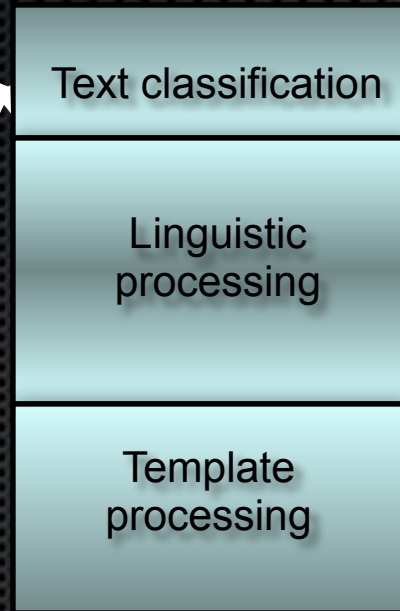
TimeIn: \_\_\_\_\_

TimeOut: \_\_\_\_\_



documents

**Dr. Hermann Wirth**, bisheriger **Leiter** der **Musikhochschule München**, verabschiedete sich heute aus dem Amt. Der 65jährige tritt seinen wohlverdienten Ruhestand an. Als seine Nachfolgerin wurde **Sabine Klinger** benannt. Ebenfalls neu besetzt wurde die Stelle des Musikdirektors. Annelie Häfner folgt Christian Meindl nach.



ManagementSuccession

PersonIn: Klinger

PersonOut: Wirth

Position: Leiter

Organisation: Musikhochschule München

TimeIn: \_\_\_\_\_

TimeOut: 3.4.2002



# Example: LIEP (S. Huffman, 1995)

<PNG> Sue Smith </PNG>, 39, of Menlo Park, was appointed <TNG> president </TNG> of <CNG> Foo Inc. </CNG>

n\_was\_named\_t\_by\_c:

noun-group(PNG, head(isa(person-name))),

noun-group(TNG, head(isa(title))),

noun-group(CNG, head(isa(company-name))),

verb-group(VG, type(passive), head(named or elected or appointed)),

prep(PREP, head(of or at or by)),

subject(PNG, VG), object(VG, TNG), post\_nominal\_prep(TNG, PREP), prep\_obj  
(PREP, CNG)

⇒ management\_appointment(M, person(PNG), title(TNG), company(CNG))



# Major IE tasks

- ✦ Named Entity task (NE)
- ✦ Template Element task (TE)
- ✦ Template Relation task (TR)
- ✦ Scenario Template task (ST)
- ✦ Co-reference task (CO)





# Named Entity Task (NE)

Mark into the text each string that represents a person, organization, or location name, or a date or time, or a currency or percentage figure (this classification of NEs reflects the standard types of NE applied in IE).

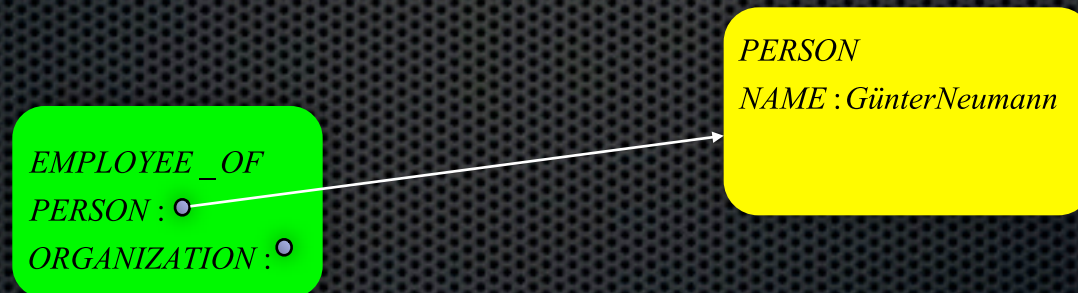
*PERSON*

*NAME : GünterNeumann*



# Template Element Task (TE)

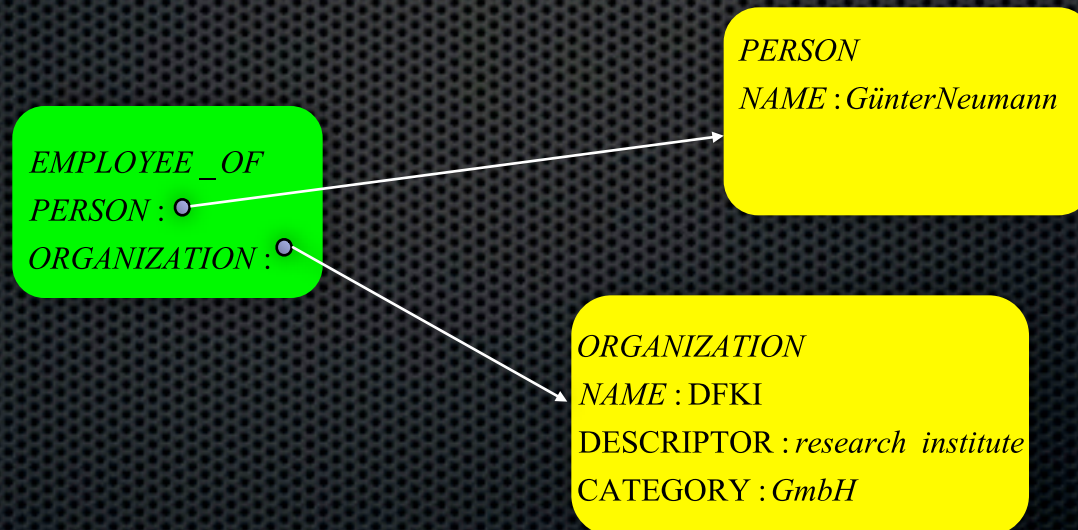
Extract basic information related to organization, person, and artifact entities, drawing evidence from everywhere in the text (also known as slot filler task; it is basically a binary relation of form „attribute of x has value y“)





# Template Relation task (TR)

Extract relational information on employee\_of, manufacture\_of, location\_of relations etc. (TR expresses domain-independent relationships between entities identified by TE)



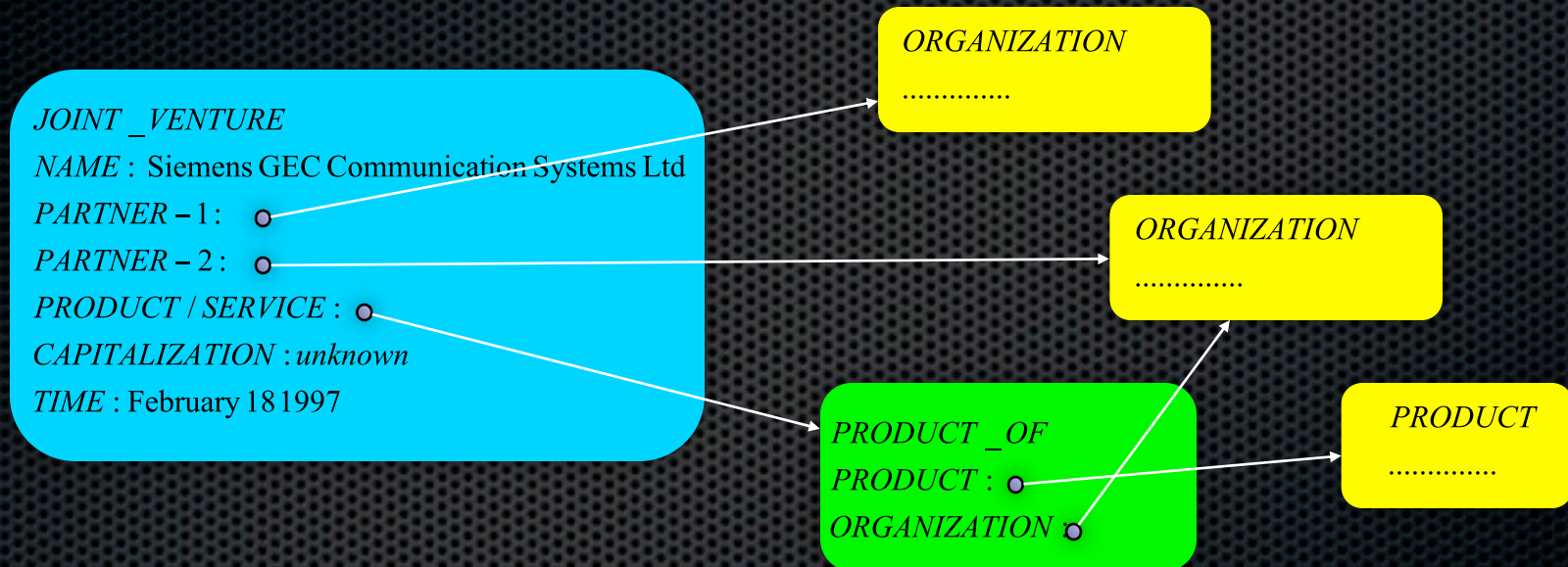


# Scenario Template task (ST)

Extract pre-specified event information and relate the event information to particular organization, person, or artifact entities (ST identifies domain and task specific entities and relations)



# ST example





# Coreference task (CO)

Capture information on co-referring expressions, i.e. all mentions of a given entity, including those marked in NE and TE (Nouns, Noun phrases, Pronouns)



Prince

He

The Artist Formerly Known As Prince



# An Example

The shiny red rocket was fired on Tuesday. It is the brainchild of Dr. Big Head. Dr. Head is a staff scientist at We Build Rockets Inc.

- NE: *red rocket, Tuesday, Dr. Big Head, Dr. Head, and We Build Rockets Inc.*
- CO: *it* refers to the rocket; *Dr. Head* and *Dr. Big Head* are the same
- TE: the rocket is *shiny red* and Head's *brainchild*
- TR: Dr. Head *works for* We Build Rockets Inc.
- ST: a *rocket launching event* occurred with the various participants.



# Scoring templates

- ✧ Templates are compared on a slot-by-slot basis
  - ✧ Correct:  $\text{response} = \text{key}$
  - ✧ Partial:  $\text{response} \approx \text{key}$
  - ✧ Incorrect:  $\text{response} \neq \text{key}$
  - ✧ Spurious: key is blank
    - ✧  $\text{overgen} = \text{spurious} / \text{actual}$
  - ✧ Missing: response is blank



# Evaluation Metrics

- ✦ Precision and recall:

- ✦ Precision: correct answers/answers produced
- ✦ Recall: correct answers/total possible answers

- ✦ F-measure

- ✦ Where  $\beta$  is a parameter representing relative importance of P & R:

$$F = \frac{(\beta^2 + 1)PR}{(\beta^2 P + R)}$$

- ✦ E.g.,  $\beta=1$ , then P&R equal weight,  $\beta=0$ , then only P

- ✦ Current State-of-Art: **F=.60 barrier**



# Maximum Results Reported in MUC-7 (Message Understanding Conference, 2001)

Measure\Task	NE	CO	TE	TR	ST
Recall	92	56	86	67	42
Precision	95	69	87	86	65

Human on NE task	F	R	P
Annotator 1	98.6	98	98
Annotator 2	96.9	96	98

Human on ST task: ~ 80 % F

This is mainly for hand-written systems



# Jerry Hobbs: Why the 60% Barrier

1. Merging problems accounted for 60% of our errors.
2. Entity recognition performance is 90%;  
event recognition requires recognizing ~4 entities;  
 $.9^4 = .6$
3. The distribution of problems has a very long tail.
4. 60% is what the text wears on its sleeve; the rest is implicit and requires inference and world knowledge.

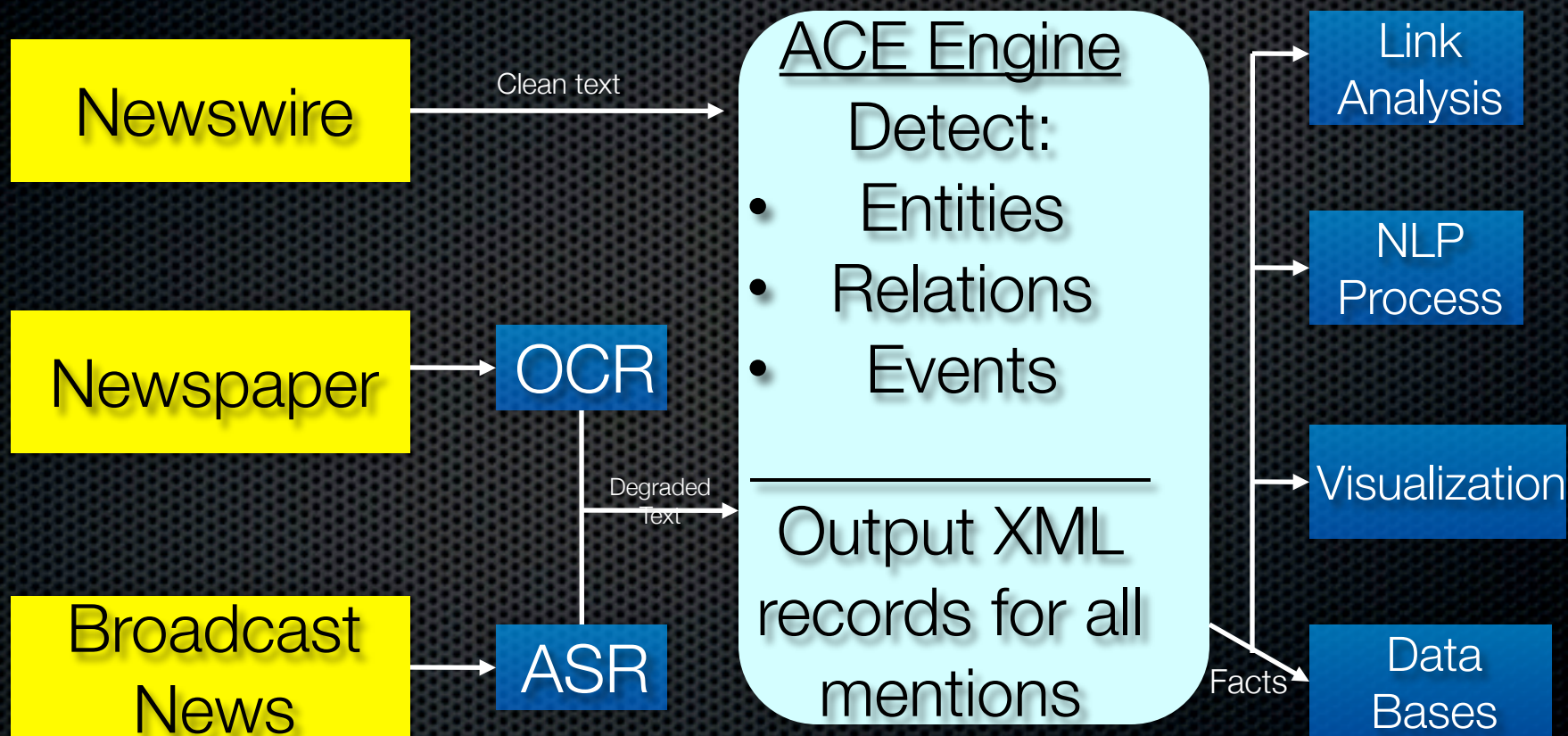


# Automatic Content Extraction - ACE (cf. Appelt, 2003)

- Develop **core information extraction technology** by focusing on extracting specific semantic entities and relations over a very wide range of texts.
- Corpora: Newswire and broadcast transcripts, but **broad range of topics and genres**.
  - Third person reports
  - Interviews
  - Editorials
  - Topics: foreign relations, significant events, human interest, sports, weather
- **Discourage** highly domain- and genre-dependent solutions



# The Technical Approach





# Objectives and Performance Goals of ACE

- Objectives
  - Extract Info from Texts of Varying Quality
  - Detect Unique entities, relations, events
    - Find **all mentions** within documents
    - Collect all mentions by object
  - **Track entities** within & across documents
  - Output XML for follow-on processes
- Performance Goals
  - Extract **95%** of the value in document



# Example for entities & their mentions

[COLOGNE, [Germany]] (AP) [A [Chilean] exile] has filed a complaint against former [Chilean] dictator Gen. Augusto Pinochet accusing [him] of responsibility for [her] arrest and torture in [Chile] in 1973, [prosecutors] said Tuesday.

[The woman, [a Chilean] who has since gained [German] citizenship] accused [Pinochet] of depriving [her] of personal liberty and causing bodily harm during [her] arrest and torture.

Person

Organization

Geopolitical Entity



# Core Mission: Information Gathering

- ✦ Information content is main interest of human language text
  - ✦ **Semantics** drives information gathering
  - ✦ Syntax is the vehicle for organizing the information
- ✦ ACE systems provide NL understanding
  - ✦ **Detect** each entity, relation, and event of specific type
  - ✦ **Recognize** all mentions of entities, relations & events
  - ✦ **Resolve** all mentions to the proper entity, relation, or event
- ✦ Convert information in human language into structured data
  - ✦ Extract semantics of communication
  - ✦ Output in ACE program format
- ✦ Structured data supports real world modeling & analysis



# Components of a Semantic Model

- Entities - Individuals in the world *that are mentioned in a text*
  - Simple entities: singular objects
  - Collective entities: sets of objects of the same type *where the set is explicitly mentioned in the text*
- Attributes - Timeless unary properties of entities (e.g. Name)
- Temporal points and intervals
- Relations - Properties that hold of two entities over a time interval
- Events - A particular kind of relation among entities implying a change in relation state at the end of the time interval.



# Semantic Analysis: Relating Language to the Model

- Linguistic Mention
  - A particular linguistic phrase
  - Denotes a particular entity, relation, or event
    - A noun phrase, name, or possessive pronoun
    - A verb, nominalization, compound nominal, or other linguistic construct relating other linguistic mentions
- Linguistic Entity
  - **Equivalence class** of mentions with same meaning
    - **Co-referring** noun phrases
  - Relations and events derived from different mentions, but conveying the **same meaning**



# Choosing an Ontology for IE Semantics

- Ordinary native speakers should be able to annotate text with minimal training.
- People should have well-developed intuitions about type classification
  - Is a “museum” an organization or facility?
- People should have well-developed intuitions about entity coreference
  - “Peace in the Middle East”
- Entities should be extensional, not abstract, generic, counterfactual, or fictional



# Relations

- Relations hold between two entities over a time interval.
- Relations may be “timeless” or temporal interval is not specified
- Relations have inertia, i. e. they don’t change unless a relevant event happens.



# Explicit and Implicit Relations

- ✦ Many relations are true in the world. Reasonable knowledge bases used by extraction systems will include many of these relations. Semantic analysis requires focusing on certain ones that are directly motivated by the text.
- ✦ Example:
  - ✦ Baltimore is in Maryland is in United States.
  - ✦ “Baltimore, MD”
  - ✦ Text mentions Baltimore and United States. Is there a relation between Baltimore and United States?



# Explicit Relations

- ✦ Explicit relations are expressed by certain surface linguistic forms
  - ✦ Copular predication - Clinton **was** the president.
  - ✦ Prepositional Phrase - The CEO **of** Microsoft...
  - ✦ Prenominal modification - The **American** envoy...
  - ✦ Possessive - Microsoft'**s** chief scientist...
  - ✦ SVO relations - Clinton arrived in Tel Aviv...
  - ✦ Nominalizations - **Anan's visit** to Baghdad...
  - ✦ Apposition - [Tony Blair, [Britain's prime minister]...]



# Textual Analysis Conference (since 2009)

## - Knowledge Base Population Track



### Motivation

- **IE & QA technologies have been studied in isolation**
  - Not focused on discovery of information for inclusion in an existing knowledge base
  - No consideration of novelty, contradiction
- **Issues when filling in a KB**
  - Accurate extraction of facts
  - Global resolution of entities
  - Maintaining provenance of asserted facts
  - Avoiding contradiction / detection of novel information
  - Temporal qualification of assertions
  - Leveraging existing KB to assist with extraction
  - Scalability



# Textual Analysis Conference (since 2009)

## - Knowledge Base Population Track



### Comparison to ACE & TREC-QA

---

- **Corpus vs. document focus**
  - ACE: component tasks (NER, relation extraction) for a set of isolated documents
  - KBP: learn facts from a corpus. Repetition not very important. Asserting wrong information is bad.
- **Context**
  - In KBP, there is a reference knowledge base, so avoiding redundancy and detecting contradiction are important
  - In KBP slots are fixed and targets change. In TREC QA, the targets dictated which questions were asked.
- **Knowing when you don't know**
  - TREC QA had a small percentage of NIL questions (4-10%)



# Textual Analysis Conference (since 2009)

## - Knowledge Base Population Track



### KBP Snapshot

---

- **Track structure**
  - NIST – overall organization, infrastructure, evaluation
  - LDC – develop and distribute data resources, target selection, human assessments
- **Datasets**
  - LDC produced 1.3M English newswire collection
  - Reference KB populated with semi-structured facts obtained from English Wikipedia (Oct '08 dump)
    - 200k people, 200k GPEs, 60k orgs, 300+k misc/non-entities
- **Two tasks**
  - **Entity Linking** - Grounding entity mentions in documents to KB entries
  - **Slot Filling** - Learning attributes about target entities



# Textual Analysis Conference (since 2009)

- Knowledge Base Population Track



## Sample KB Entry

```
<entity wiki_title="Michael_Phelps"
  type="PER"
  id="E0318992"
  name="Michael Phelps">
<facts class="Infobox Swimmer">
<fact name="swimmername">Michael Phelps</fact>
<fact name="fullname">Michael Fred Phelps</fact>
<fact name="nicknames">The Baltimore Bullet</fact>
<fact name="nationality">United States</fact>
<fact name="strokes">Butterfly, Individual Medley, Freestyle, Backst
<fact name="club">Club Wolverine, University of Michigan</fact>
<fact name="birthdate">June 30, 1985 (1985-06-30) (age 23)</fact>
<fact name="birthplace">Baltimore, Maryland, United States</fact>
<fact name="height">6 ft 4 in (1.93 m)</fact>
<fact name="weight">200 pounds (91 kg)</fact>
</facts>
```

```
<wiki_text><![CDATA[Michael Phelps
Michael Fred Phelps (born June 30, 1985) is an American swimmer. H
Olympic gold medals, the most by any Olympian. As of August 2008,
world records in swimming. Phelps holds the record for the most gol
single Olympics with the eight golds he won at the 2008 Olympic Gar
```

Michael Phelps



Michael Phelps at the 2008 Beijing Olympics

### Personal information

Full name:	Michael Fred Phelps
Nickname(s):	The Baltimore Bullet <sup>1</sup>
Nationality:	United States
Stroke(s):	Butterfly, Individual Medley, Freestyle, Backstroke
Club:	Club Wolverine, University of Michigan
Date of birth:	June 30, 1985 (age 23)
Place of birth:	Baltimore, Maryland, United States
Height:	6 ft 4 in (1.93 m)
Weight:	200 pounds (91 kg)

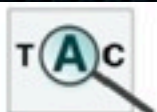
### Medal record

[show](#)



# Textual Analysis Conference (since 2009)

## - Knowledge Base Population Track



### Most Frequent KB Classes

95142 settlement	8353 ort in deutschland	5222 lake
72992 album	8061 university	4913 television episode
34659 film	7675 airport	4636 school
32464 musical artist	7492 military person	4426 commune de france
23138 actor	7270 road	4265 aircraft
21195 single	7185 indian jurisdiction	4229 ice hockey player
16765 company	7123 cityit	3918 german location
15644 book	6143 australian place	3234 nflactive
14567 football biography	6131 mountain	3168 disease
14121 person	5957 military conflict	3070 politician
12646 radio station	5952 military unit	3036 u.s. county
12514 nrhp	5937 city	2956 station
12324 vg	5630 software	2950 automobile
11813 planet	5501 mlb retired	2933 officeholder
10818 uk place	5397 writer	2833 broadcast
10113 television	5349 scientist	2728 swiss town

PER

ORG

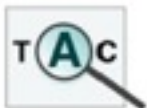
GPE

OTHER



# Textual Analysis Conference (since 2009)

## - Knowledge Base Population Track



### Entity Linking Task

#### John Williams

Richard Kaufman goes a long way back with **John Williams**. Trained as a classical violinist, Californian Kaufman started doing session work in the Hollywood studios in the 1970s. One of his movies was Jaws, with **Williams** conducting his score in recording sessions in 1975...

#### Michael Phelps

Debbie Phelps, the mother of swimming star **Michael Phelps**, who won a record eight gold medals in Beijing, is the author of a new memoir, ...

**Michael Phelps** is the scientist most often identified as the inventor of PET, a technique that permits the imaging of biological processes in the organ systems of living individuals. **Phelps** has ...



John Williams	author	1922-1994
J. Lloyd Williams	botanist	1854-1945
John Williams	politician	1955-
John J. Williams	US Senator	1904-1988
John Williams	Archbishop	1582-1650
<b>John Williams</b>	<b>composer</b>	<b>1932-</b>
Jonathan Williams	poet	1929-

Michael Phelps	swimmer	1985-
Michael Phelps	biophysicist	1939-

Identify matching entry, or determine that entity is missing from KB



# Textual Analysis Conference (since 2009)

## - Knowledge Base Population Track



### Slot Filling Task

**Target: EPA**  
(plus 1 document)

**Generic Entity Classes**  
Person, Organization, GPE

Environmental Protection Agency



Agency overview

Employees	17,964 (2005)
Annual budget	\$7.3 billion (2007)
Agency executive	Lisa P. Jackson, Administrator

**Missing information to mine from text:**

- Date formed: **12/2/1970**
- Website: **<http://www.epa.gov/>**
- Headquarters: **Washington, DC**
- Nicknames: **EPA, USEPA**
- Type: **federal agency**
- Address: **1200 Pennsylvania Avenue NW**

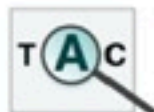
**Optional: Also want to link some learned values within the KB:**

- Headquarters: **Washington, DC (kbid: 735)**

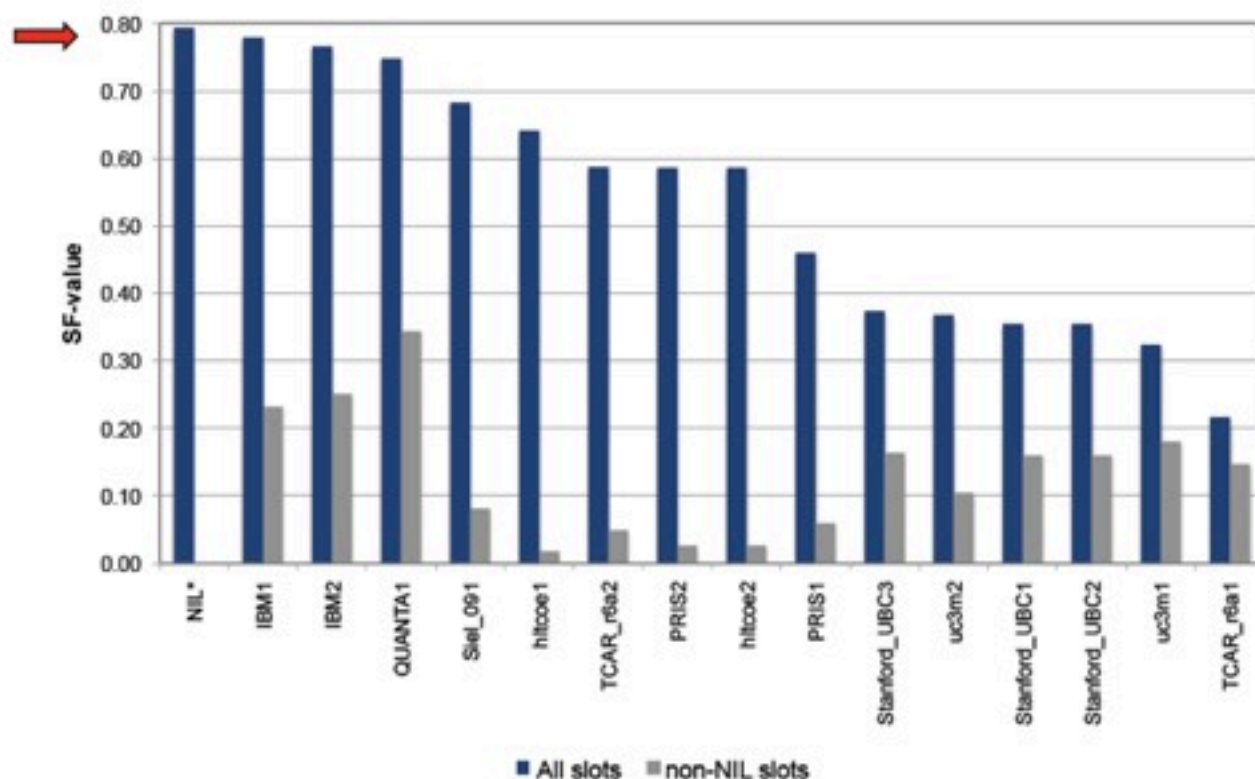


# Textual Analysis Conference (since 2009)

- Knowledge Base Population Track



## SF Results





# Two Approaches to Building Extraction Systems

- ✦ Knowledge engineering approach
  - ✦ Grammars are constructed by hand
  - ✦ Domain patterns are discovered by a human expert through introspection and inspection of a corpus
  - ✦ Much laborious tuning and „hill climbing“
- ✦ Automatically Trainable Systems
  - ✦ Use statistical methods when possible
  - ✦ Learn rules from annotated corpora
  - ✦ Learn rules from interaction with user



# Knowledge Engineering

## ✦ Advantages

- ✦ With skill and experience, good performing systems are conceptually not hard to develop
- ✦ The best performing systems have been hand crafted (still true for scenario patterns)

## ✦ Disadvantages

- ✦ Very laborious development process
- ✦ Domain adaptation might require re-configuration
- ✦ Needs experts which have both, linguistic & domain expertise



# Hand-Coded Methods

- ✦ Easy to construct in many cases
  - ✦ e.g., to recognize prices, phone numbers, conference names, etc.
- ✦ Easier to debug & maintain
  - ✦ especially if written in a “high-level” language (as is usually the case)
  - ✦ e.g.,

```
ContactPattern ← RegularExpression(Email.body, "can be reached at")  
  
PersonPhone ← Precedes(Person  
                        Precedes(ContactPattern, Phone, D),  
                        D)
```

- ✦ Easier to incorporate / reuse domain knowledge
- ✦ Can be quite labor intensive to write



# Learning-Based Methods

- ✦ Can work well when training data is easy to construct and is plentiful
- ✦ Can capture complex patterns that are hard to encode with hand-crafted rules
  - ✦ e.g., determine whether a review is positive or negative
  - ✦ extract long complex gene names

*The **human T cell leukemia lymphotropic virus type 1 Tax protein** represses MyoD-dependent transcription by inhibiting MyoD-binding to the KIX domain of p300.*“

- ✦ Can be labor intensive to construct training data
  - ✦ not sure how much training data is sufficient



# Trainable Systems

## ✦ Advantages

- ✦ Domain portability is relatively straightforward
- ✦ System expertise is not required for customization
- ✦ Data driven rule acquisition ensures full coverage of examples

Possible solutions here are

- on-demand IE
- dynamic interactive IE
- ad-hoc IE

## ✦ Disadvantages

- ✦ Training data may not exist, and maybe very expensive to acquire
- ✦ Large volume of training data may be required
- ✦ Changes to specifications may require re-annotation of large quantities of training data



# What works best?

- Use rule-based approach when
  - Resources (e.g., lexicons, lists) are available
  - Rule writers are available
  - Training data scarce or expensive to obtain
  - Extraction specs likely to change
  - Highest possible performance is critical
- Use trainable approach when
  - Resources unavailable
  - No skilled rule writers are available
  - Training data is cheap and plentiful
  - Good performance is adequate for the task

This is still the main approach in industrial applications, where precision counts.



# Architecture of Extraction Systems

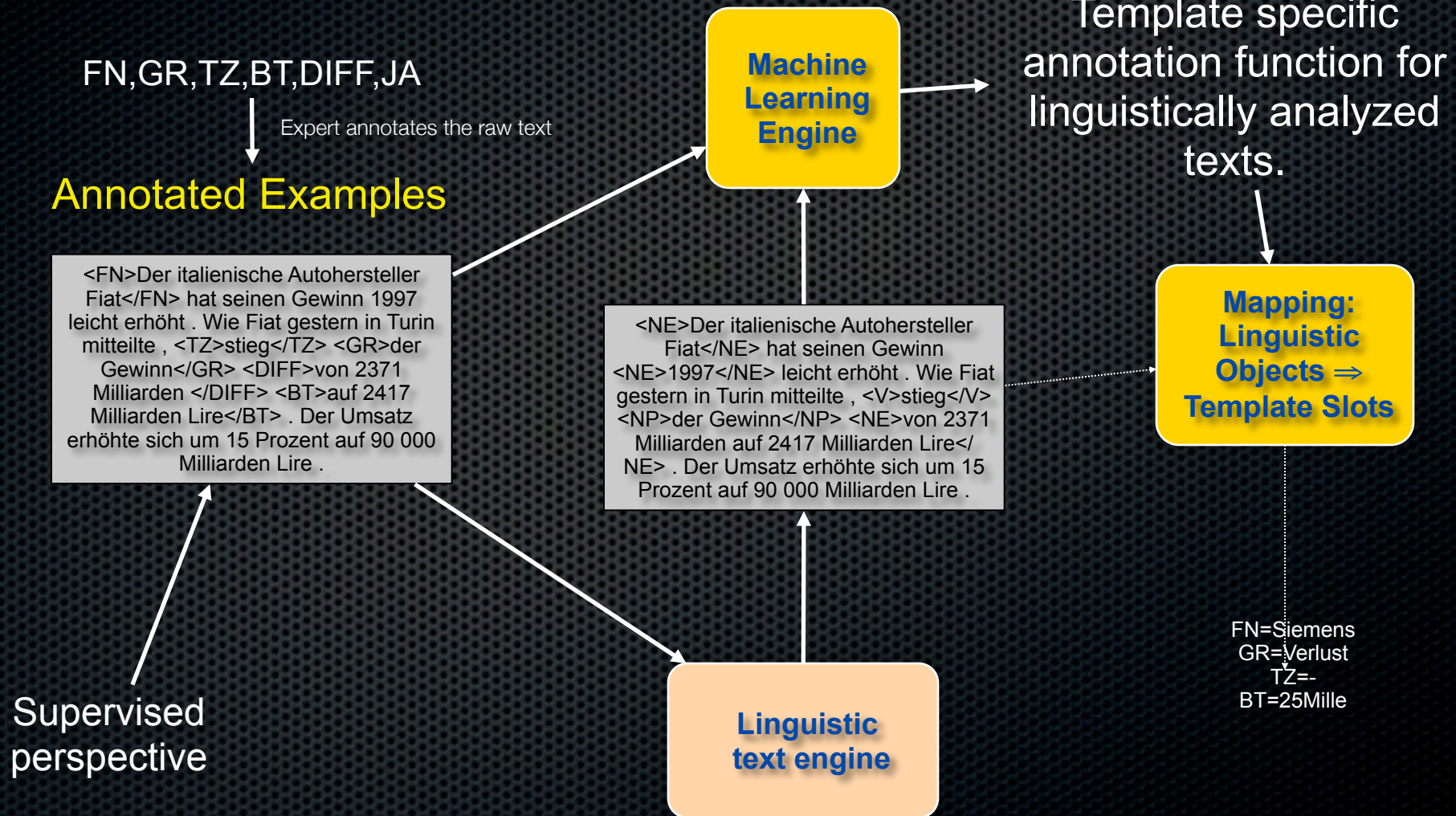
- Domain-independent NL tools necessary
  - Main task: define scope of linguistic features
  - Major issue: robustness & efficiency
- Clean interface between domain-independent tools and domain-dependent
  - Domain modeling
  - Main Task: disambiguation
  - Easy adaptation of NL tools



# IE and Machine Learning

Input: Template specification  
(e.g., company turnover/revenue)

Output:  
Template specific  
annotation function for  
linguistically analyzed  
texts.

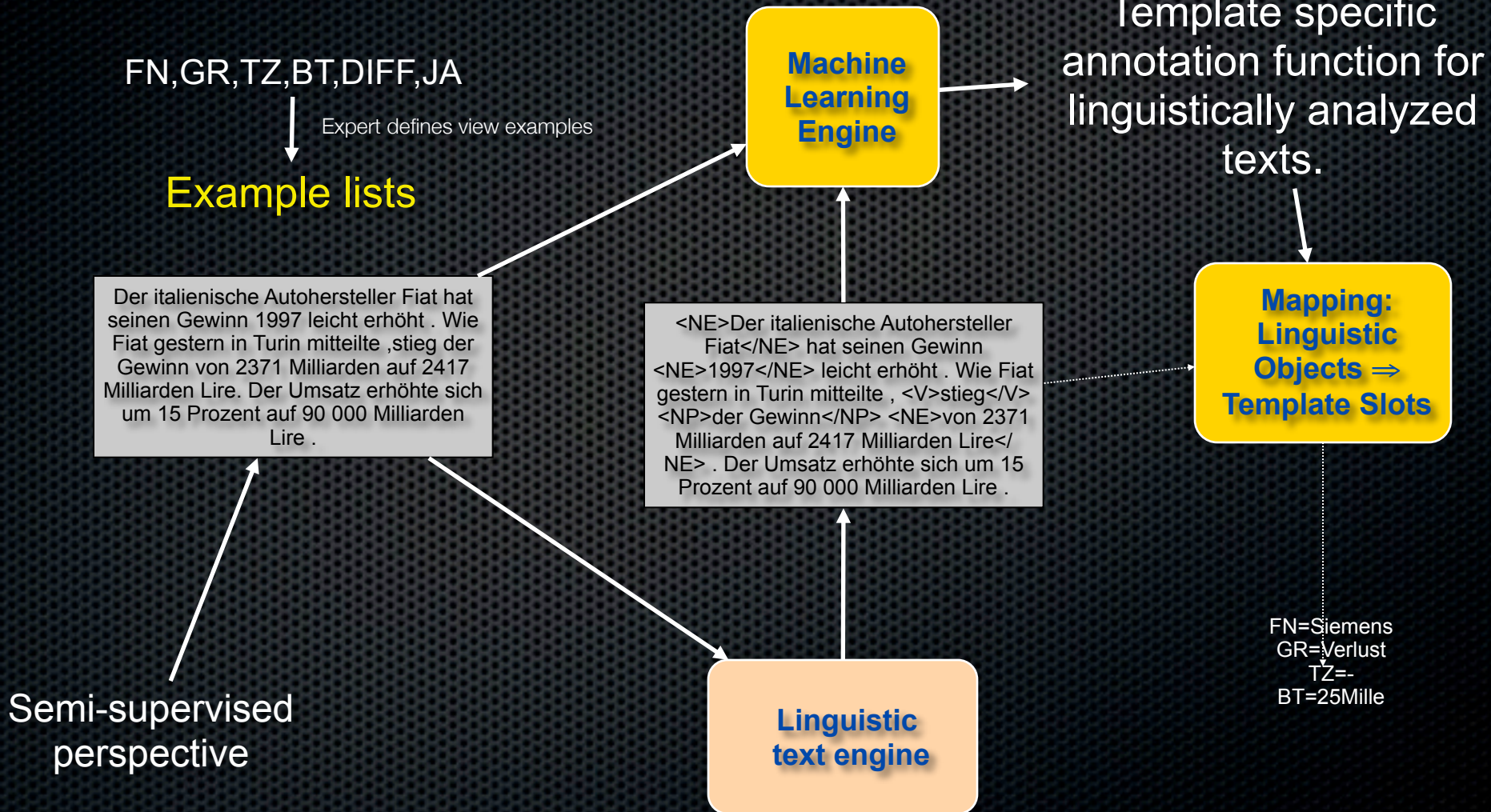




# IE and Machine Learning

Input: Template specification  
(e.g., company turnover/revenue)

Output:  
Template specific  
annotation function for  
linguistically analyzed  
texts.



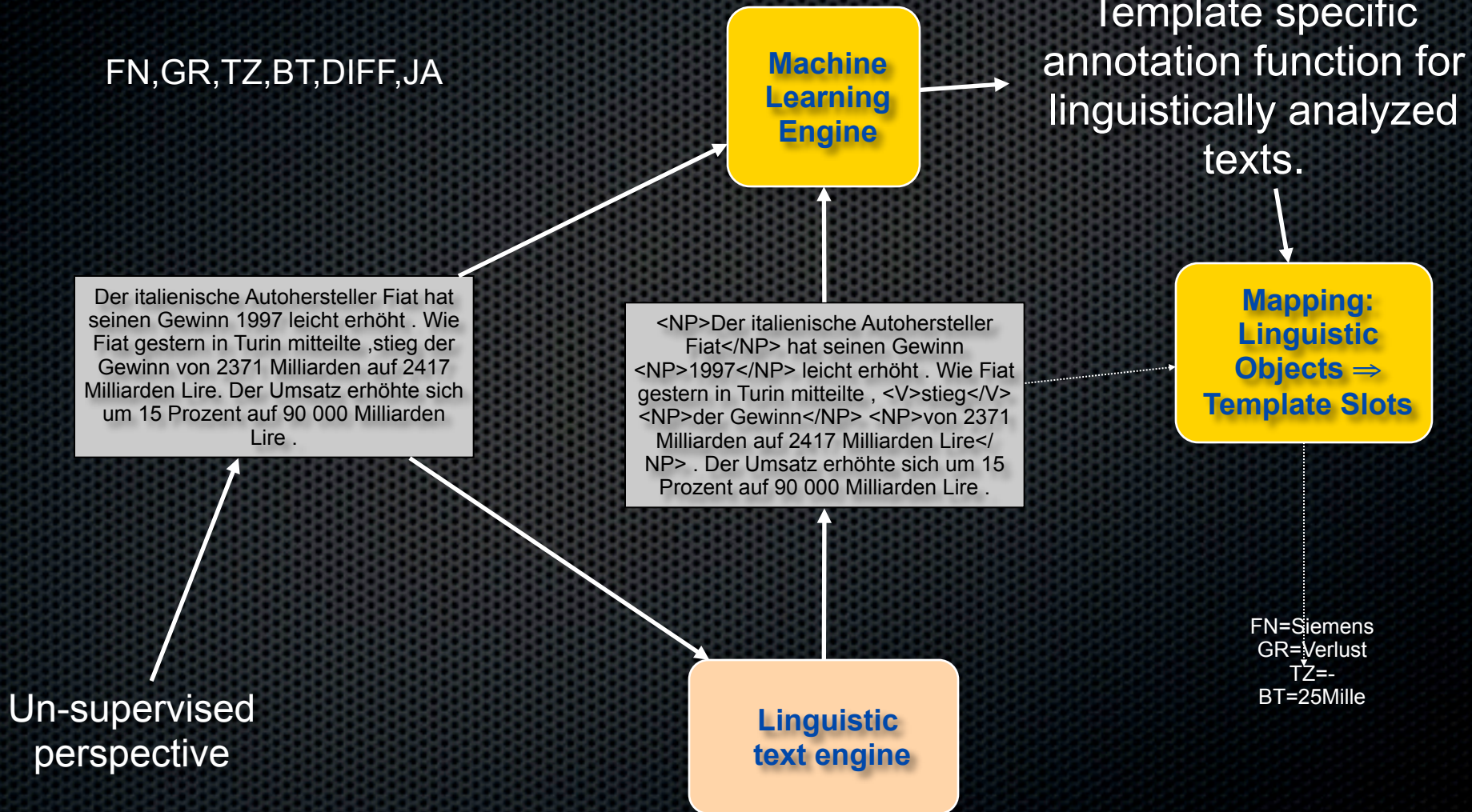


# IE and Machine Learning

Input: Template specification  
(e.g., company turnover/revenue)

FN,GR,TZ,BT,DIFF,JA

Output:  
Template specific  
annotation function for  
linguistically analyzed  
texts.





# In any case, major subtasks

- ✦ Linguistic feature extraction
- ✦ Recognition and extraction of Named Entities
- ✦ Recognition and extraction of Relations
- ✦ Recognition and extraction of Events



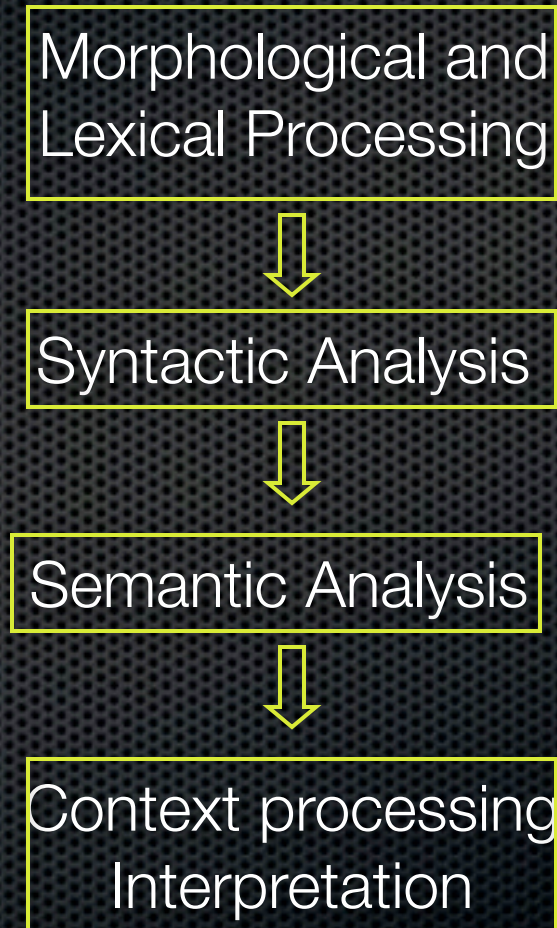
# About the Linguistic Core Engine

- ✦ From part-of-speech, to morphology, to phrase recognition, to full sentence structure to logical formulas.
- ✦ What linguistic features are needed or useful depends on the application task and strategy.



# General Framework of NLP

John runs.





# General Framework of NLP

John runs.

John run+s.

P-N    V    3-pre  
         N    plu

Morphological and  
Lexical Processing



Syntactic Analysis



Semantic Analysis



Context processing  
Interpretation



# General Framework of NLP

John runs.

John run+s.

P-N    V    3-pre  
      N    plu

Morphological and  
Lexical Processing



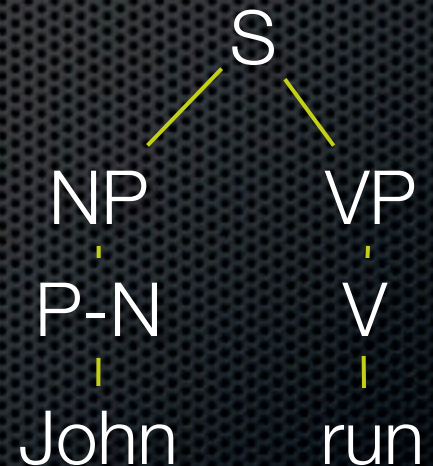
Syntactic Analysis



Semantic Analysis



Context processing  
Interpretation





# General Framework of NLP

John runs.

John run+s.

P-N    V    3-pre  
         N    plu

Pred: RUN  
Agent: John

Morphological and  
Lexical Processing



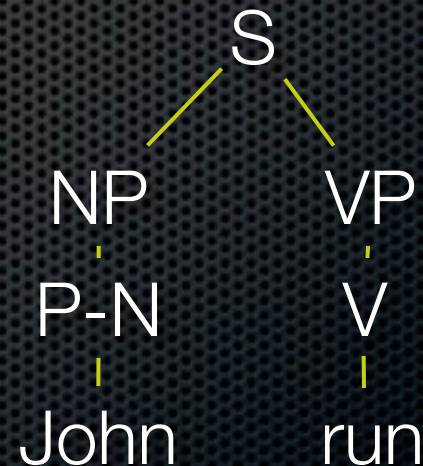
Syntactic Analysis



Semantic Analysis



Context processing  
Interpretation





# General Framework of NLP

John runs.

John run+s.

P-N   V   3-pre  
      N   plu

Pred: RUN  
Agent: John

John is a student.  
He runs.

Morphological and  
Lexical Processing



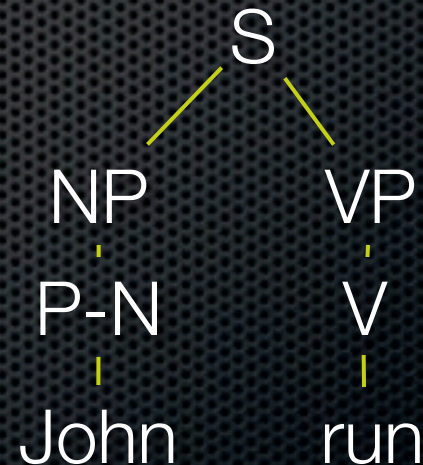
Syntactic Analysis



Semantic Analysis



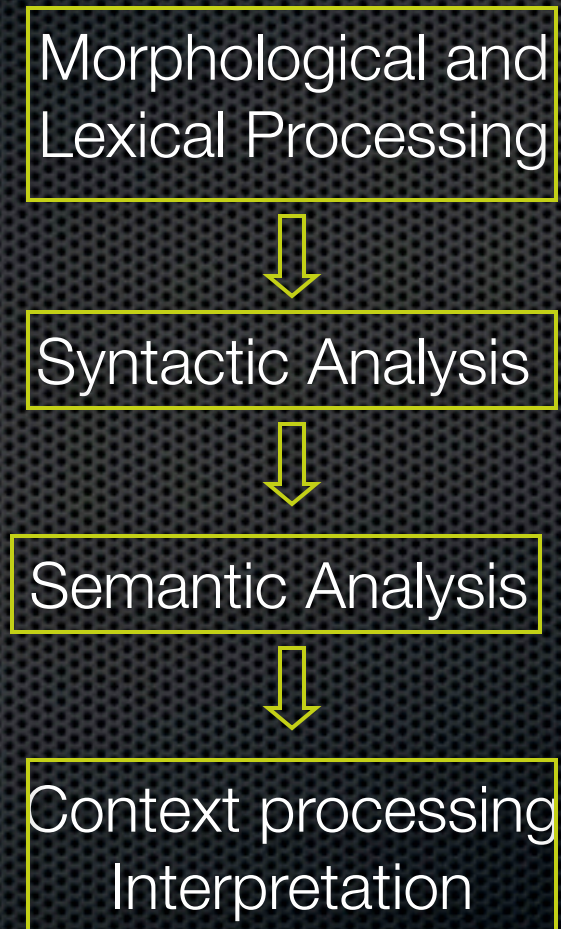
Context processing  
Interpretation





# *Difficulties of NLP*

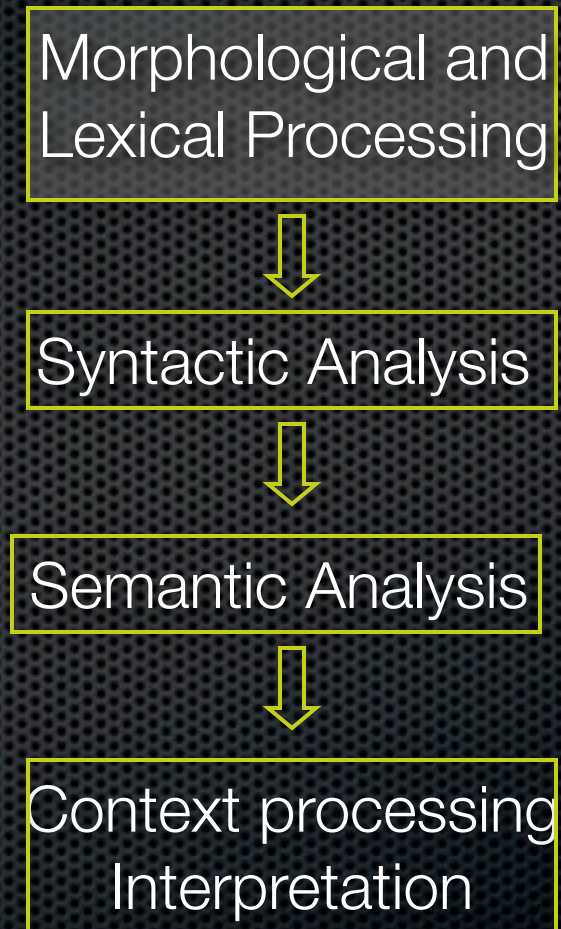
(1) Robustness: General Framework of NLP  
Incomplete Knowledge





# *Difficulties of NLP*

(1) Robustness: General Framework of NLP  
Incomplete Knowledge



Incomplete Lexicons  
Open class words  
Terms  
Term recognition  
Named Entities  
Company names  
Locations  
Numerical expressions



# *Difficulties of NLP*

(1) Robustness: General Framework of NLP  
Incomplete Knowledge

Incomplete Grammar  
Syntactic Coverage  
Domain Specific  
Ungrammatical

Morphological and  
Lexical Processing



Syntactic Analysis



Semantic Analysis

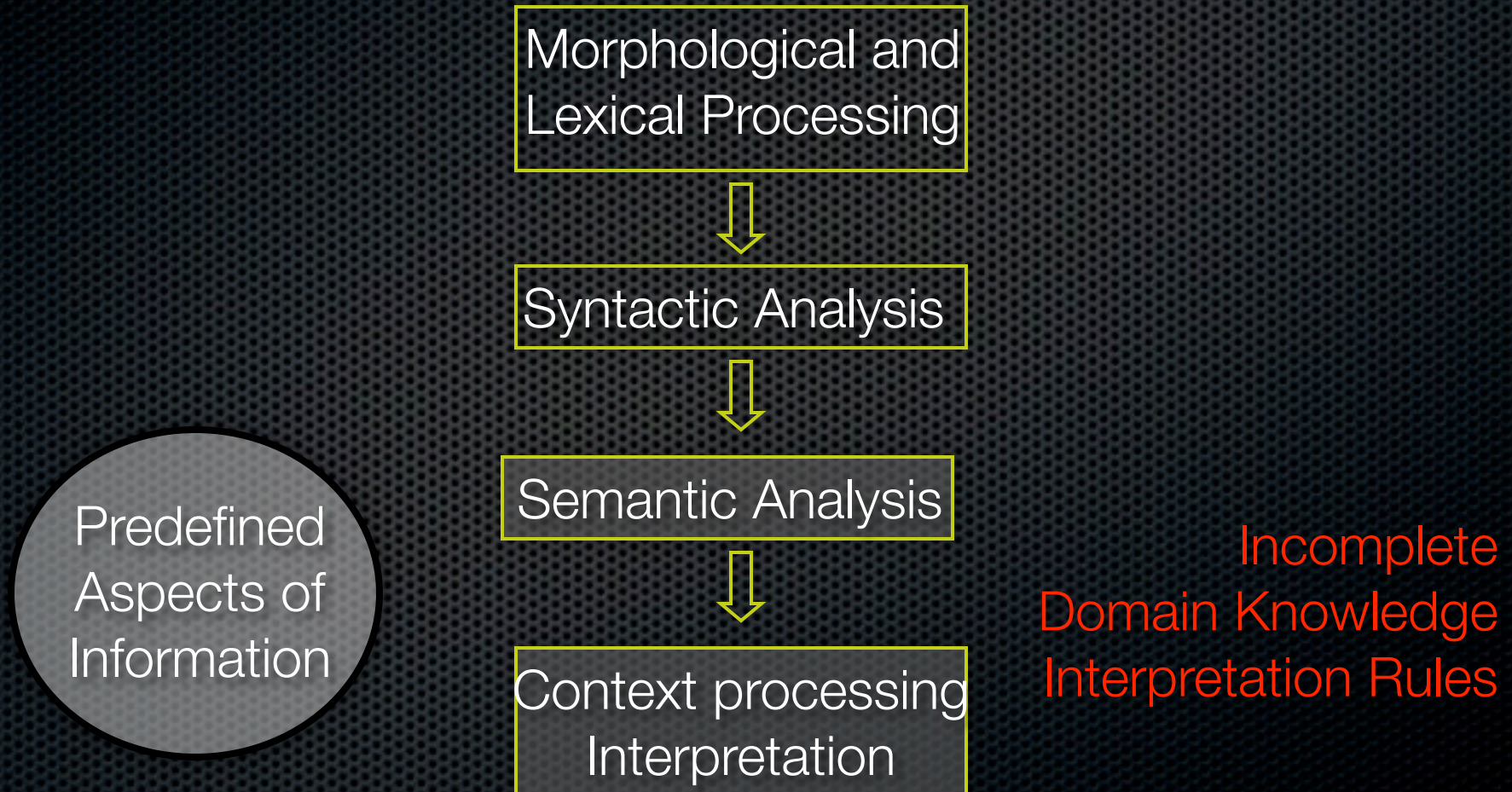


Context processing  
Interpretation



# *Difficulties of NLP*

(1) Robustness: General Framework of NLP  
Incomplete Knowledge

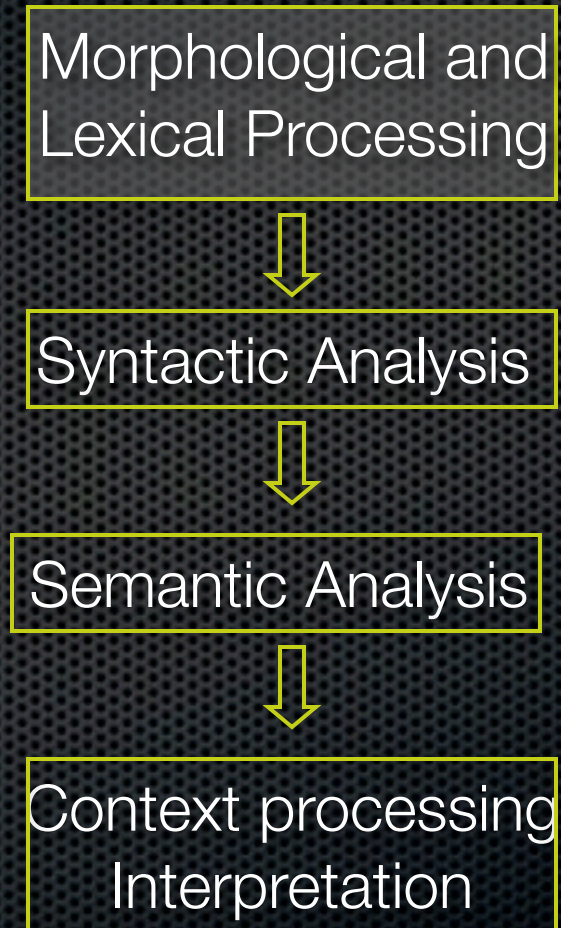




# *Difficulties of NLP*

(1) Robustness: General Framework of NLP  
Incomplete Knowledge

(2) Ambiguities:  
Combinatorial  
Explosion



Most words in English  
are ambiguous in terms  
of their part of speeches  
runs: v/3pre, n/plu  
clubs: v/3pre, n/plu  
and two meanings



# *Difficulties of NLP*

(1) Robustness: General Framework of NLP  
Incomplete Knowledge

(2) Ambiguities:  
Combinatorial  
Explosion

*Combinatorial  
Explosion*

Morphological and  
Lexical Processing



Syntactic Analysis



Semantic Analysis



Context processing  
Interpretation

Structural Ambiguities

Predicate-argument  
Ambiguities

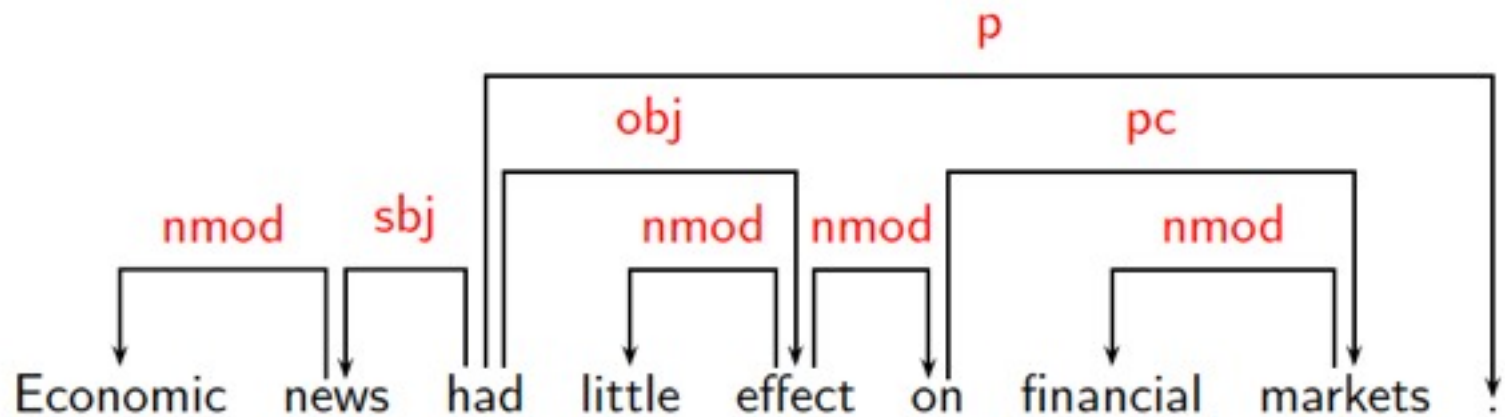


# Dependency Parsing

- The basic idea:
  - Syntactic structure consists of lexical items , linked by binary asymmetric relations called dependencies
  - The dependency structure of a sentence is a acyclic directed graph (DAG; mostly, they are such trees)
- A dependency parsing
  - computes a dependency structure for a sentence
  - tree-bank-based dependency parsers: the rules or constraints for determining the valid structure is learned from a large set of example valid dependency structures (trees)

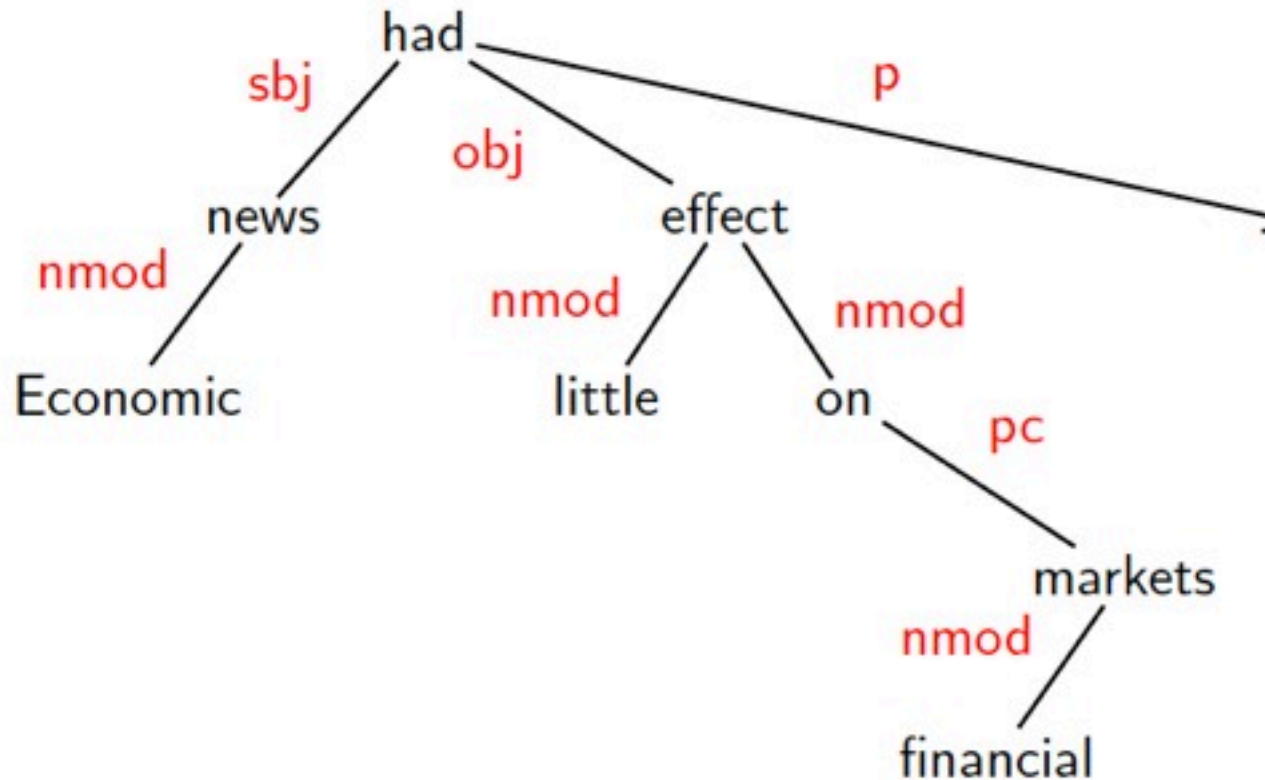


# Dependency Structure





# Notational Variants





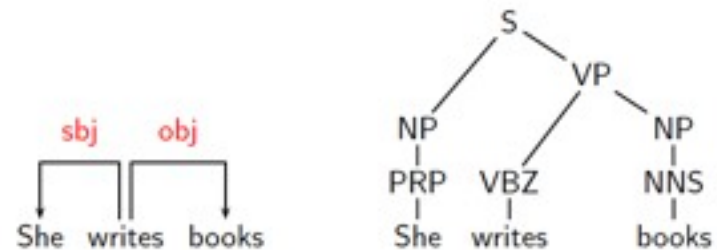
# Comparison

- Dependency structures explicitly represent
  - head-dependent relations (directed arcs ),
  - functional categories (arc labels ),
  - possibly some structural categories (parts-of-speech).
- Phrase structures explicitly represent
  - phrases (nonterminal nodes ),
  - structural categories (nonterminal labels ),
  - possibly some functional categories (grammatical functions).



# Application-oriented Advantage

- Shallow semantics: Direct encoding of predicate-argument structure
- Machine Learning:
  - feature extraction
  - Tree kernels





# IE: compromise NLP (cf. Appelt & Israel, 1999; Neumann & Piskorski, 2002)

## ✦ Task characteristic

- ✦ Lots of texts
- ✦ Dirty texts
- ✦ World knowledge needed

## • Compromise

- Finite-state models
- Robust techniques
- Domain specific processing at each stage of analysis

The bottom line:

Find the most favorable tradeoff between recall and precision for the task at hand.



# Languages Are Structural

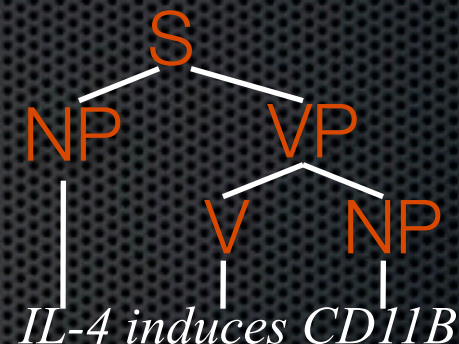
governments  
Firmen



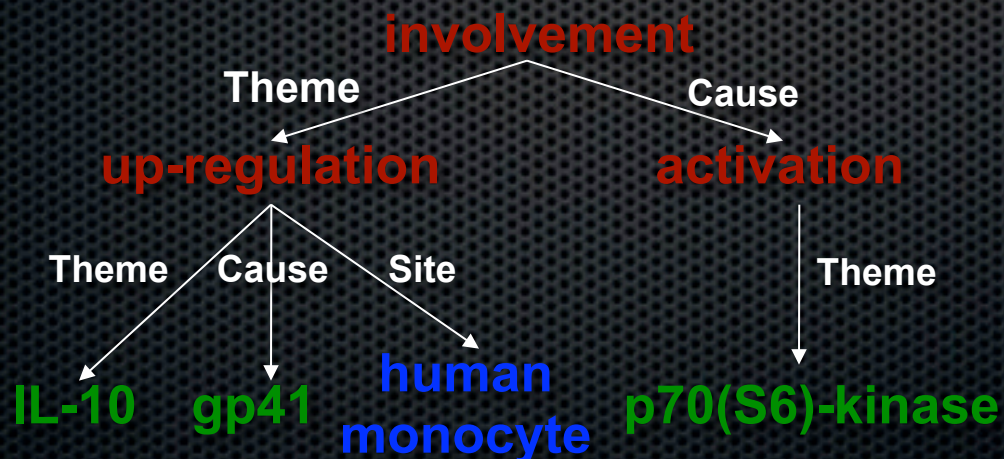
# Languages Are Structural

govern-ment-s

Firma-PL



*Involvement of p70(S6)-kinase activation in IL-10 up-regulation in human monocytes by gp41.....*



George Walker Bush was the 43<sup>rd</sup> President of the United States.

.....

Bush was the eldest son of President G. H. W. Bush and Babara Bush.

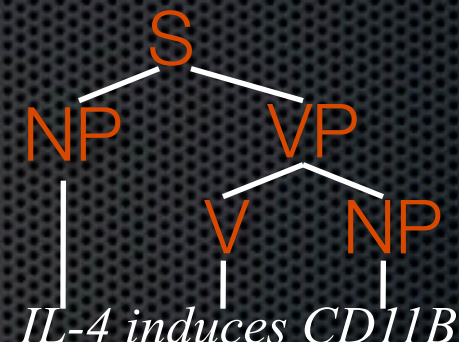
.....

In November 1977, he met Laura Welch at a barbecue.

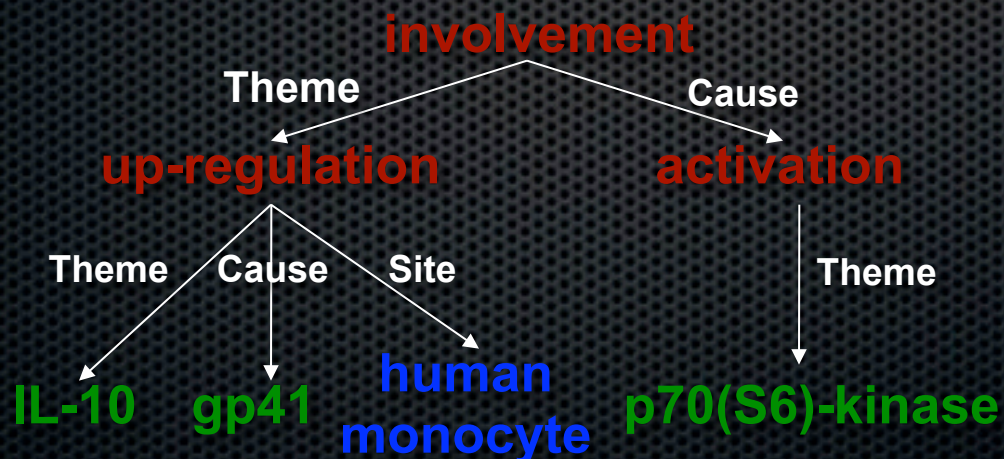


# Languages Are Structural

govern-ment-s  
l-m\$px-t-m  
(according to their families)



*Involvement of p70(S6)-kinase  
activation in IL-10 up-regulation  
in human monocytes by gp41.....*



George Walker Bush was the  
43<sup>rd</sup> President of the United  
States.

.....

Bush was the eldest son of  
President G. H. W. Bush and  
Babara Bush.

.....

In November 1977, he met  
Laura Welch at a barbecue.



# Languages Are Structural

- Objects are not just feature vectors
  - They have parts and subparts
  - Which have relations with each other
  - They can be trees, graphs, etc.
- Objects are seldom i.i.d.  
(independent and identically distributed)
  - They exhibit local and global dependencies
  - They form class hierarchies (with multiple inheritance)
  - Objects' properties depend on those of related objects
- Deeply interwoven with knowledge



# Languages Are Statistical

I saw the man with the telescope

NP

I saw the man with the telescope

NP

ADVP

I saw the man with the telescope

*Microsoft buys Powerset*

*Microsoft acquires Powerset*

*Powerset is acquired by Microsoft Corporation*

*The Redmond software giant buys Powerset*

*Microsoft's purchase of Powerset, ...*

Here in **London**, Frances Deek is a retired teacher ...

In the Israeli town ..., Karen **London** says ...

Now **London** says ...

**G. W. Bush** .....

..... **Laura Bush** .....

**Mrs. Bush** .....

**London** = PERSON or LOCATION?

Which one?



# Languages Are Statistical

- Languages are ambiguous
- Our information is always incomplete
- We need to model correlations
- Our predictions are uncertain
- Statistics provides the tools to handle this



# NL and Information Extraction

- ✦ Text ambiguity and variability is a major challenge for large-scale robust and efficient text exploration
- ✦ One has to carefully decide, which linguistic properties are actually relevant and necessary for high-coverage and high-precise extraction of information nuggets, e.g., named entities or relations between them.